

ARGONNE NATIONAL LABORATORY  
9700 South Cass Avenue  
Argonne, Illinois 60439

**DISTANCE GEOMETRY OPTIMIZATION FOR PROTEIN  
STRUCTURES**

**Jorge J. Moré and Zhijun Wu**

Mathematics and Computer Science Division

Preprint MCS-P628-1296

December 1996

This work was supported by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Computational and Technology Research, U.S. Department of Energy, under Contract W-31-109-Eng-38 and by the Argonne Director's Individual Investigator Program.



# DISTANCE GEOMETRY OPTIMIZATION FOR PROTEIN STRUCTURES

Jorge J. Moré and Zhijun Wu

## Abstract

We study the performance of the **dgsol** code for the solution of distance geometry problems with lower and upper bounds on distance constraints. The **dgsol** code uses only a sparse set of distance constraints, while other algorithms tend to work with a dense set of constraints either by imposing additional bounds or by deducing bounds from the given bounds. Our computational results show that protein structures can be determined by solving a distance geometry problem with **dgsol** and that the approach based on **dgsol** is significantly more reliable and efficient than multi-starts with an optimization code.

## 1 Introduction

Distance geometry problems for the determination of protein structures are specified by a subset  $\mathcal{S}$  of all atom pairs and by the distances between atoms  $i$  and  $j$  for  $(i, j) \in \mathcal{S}$ . In practice, lower and upper bounds on the distances are given instead of precise values. The distance geometry problem with lower and upper bounds is to find a set of positions  $x_1, \dots, x_m$  in  $\mathbb{R}^3$  such that

$$l_{i,j} \leq \|x_i - x_j\| \leq u_{i,j}, \quad (i, j) \in \mathcal{S}, \quad (1.1)$$

where  $l_{i,j}$  and  $u_{i,j}$  are lower and upper bounds on the distances, respectively. Reviews and background on the application of distance geometry problems to protein structure determination can be found in Crippen and Havel [4], Havel [11, 12], Torda and van Gunsteren [30], Kuntz, Thomason, and Oshiro [19], Brünger and Nilges [3], and Blaney and Dixon [2].

The distance geometry problem (1.1) can be formulated as a global optimization problem. The standard formulation, suggested by Crippen and Havel [4], is in terms of finding the global minimum of the function

$$f(x) = \sum_{i,j \in \mathcal{S}} p_{i,j}(x_i - x_j), \quad (1.2)$$

where the pairwise function  $p_{i,j} : \mathbb{R}^n \mapsto \mathbb{R}$  is defined by

$$p_{i,j}(x) = \min^2 \left\{ \frac{\|x\|^2 - l_{i,j}^2}{l_{i,j}^2}, 0 \right\} + \max^2 \left\{ \frac{\|x\|^2 - u_{i,j}^2}{u_{i,j}^2}, 0 \right\}. \quad (1.3)$$

---

This work was supported by the Mathematical, Information, and Computational Sciences Division subprogram of the Office of Computational and Technology Research, U.S. Department of Energy, under Contract W-31-109-Eng-38 and by the Argonne Director's Individual Investigator Program.

Clearly,  $x = \{x_1, \dots, x_m\}$  solves the distance geometry problem if and only if  $x$  is a global minimizer of  $f$  and  $f(x) = 0$ .

In practice, distance geometry problems also impose chirality constraints on some of the atoms. These constraints can be handled by adding a term to the potential function (1.2) whenever chirality constraints are imposed on the atoms  $x_i, x_j, x_k, x_l$ . The chirality constraint function (for example, Havel [11, 12]) is usually of the form

$$c_{i,j,k,l}(x) = (\text{vol}(x_i, x_j, x_k, x_l) - v_{i,j,k,l})^2,$$

where  $\text{vol}$  is the oriented volume of the four atoms and  $v_{i,j,k,l}$  is a target value. The approach in this paper can be extended to these chirality constraints, but since our aim is to develop reliable algorithms for the distance geometry problem (1.1), we consider only the potential function (1.2).

The **embed** algorithm (Crippen and Havel [4], Havel [11, 12]) and the alternating projection algorithm (Glunt, Hayden, and Raydan [7, 8]) are the most promising techniques for the solution of the distance geometry problem (1.2). For related work, see [1, 2, 3, 19]. General global optimization techniques (multi-starts with a local optimization algorithm, simulated annealing, genetic algorithms) and molecular dynamics algorithms could also be used, but they have not been shown to be suitable for distance geometry problems.

The algorithm that we propose in this paper works with the sparse set of distance constraints  $\mathcal{S}$ . In contrast, other algorithms for distance geometry tend to work with a dense set of constraints by either imposing additional bounds or by deducing bounds from the given bounds. For example, the first phase of the **embed** algorithm determines  $l_{i,j}$  and  $u_{i,j}$  by using the relationships

$$u_{i,j} = \min(u_{i,j}, u_{i,k} + u_{k,j}), \quad l_{i,j} = \max(l_{i,j}, l_{i,k} - u_{k,j}, l_{j,k} - u_{k,i}),$$

which can be deduced from the triangle inequality. Given a full set of bounds, distances  $\delta_{i,j} \in [l_{i,j}, u_{i,j}]$  are chosen, and an attempt is made to compute coordinates  $x_1, \dots, x_m$  by solving the special distance geometry problem

$$\|x_i - x_j\| = \delta_{i,j}, \quad (i, j) \in \mathcal{S}. \quad (1.4)$$

This attempt usually fails because the bounds  $\delta_{i,j}$  tend to be inconsistent, but it can be used to generate an approximate solution. As a result, the **embed** algorithm may require many trial choices of  $\delta_{i,j}$  in  $[l_{i,j}, u_{i,j}]$  before a solution to problem (1.4) is found.

Other algorithms that work with a sparse set of distance constraints do not aim to solve the distance geometry problem (1.1), but to minimize a potential energy function that incorporates distance constraints and other information to determine the protein structure. For work in this direction, see [28, 15].

In our approach, we use Gaussian smoothing to transform  $f$  into a smoother function with fewer minimizers. An optimization algorithm (the limited-memory variable-metric code **vm1m**) is then applied to the transformed function, and continuation techniques are used to trace the minimizers of the smooth function back to the original function. An immediate advantage of our approach is that the work per iteration is proportional to  $\mathcal{S}$ , which for sparse distance data should be proportional to the number of atoms  $m$ .

Gaussian smoothing was first used, by Scheraga and coworkers [25, 16, 17, 18, 26], in the diffusion equation method for protein conformation. In that application the Gaussian transform is usually evaluated by approximating the function and then transforming the approximation. On the other hand, Moré and Wu [22] showed that for distance geometry applications we can evaluate the Gaussian transform of (1.2) directly if the potential  $p_{i,j}$  is a radial function, that is, a function of the form  $p_{i,j}(x) = h_{i,j}(\|x\|)$ .

The aim of this paper is to show that continuation algorithms, based on Gaussian smoothing, can be used to develop an efficient and reliable code for the solution of the distance geometry problem (1.1). The background needed to understand our code, **dgsol**, is presented in Sections 2 and 3. Section 2 outlines the smoothing properties of the Gaussian transform, while Section 3 presents our proposal to determine the Gaussian transform by using a discrete Gauss-Hermite transform.

We present an outline of **dgsol** in Section 4. Numerical results appear in Section 5. We pay special attention to the choice of continuation parameters because this is an important and unresolved issue in the use of Gaussian smoothing. Our numerical results, based on data drawn from the PDB data bank, show that **dgsol** can be used to determine the structure of protein fragments with up to 200 atoms.

We emphasize that the determination of protein structures from distance data requires appropriate data and an algorithm to determine solutions to (1.1). The issue of what distance data is needed has been addressed in several recent papers [15, 28, 1, 20]. In this paper we do not address this issue, but concentrate on showing that **dgsol** can be used to obtain solutions to the distance geometry problem (1.1) for a wide range of distance data. To our knowledge, no other algorithm can make this claim. We plan to conduct additional testing with larger protein fragments and more realistic distance constraints.

## 2 Global Smoothing

An appealing idea for finding the global minimizer of a function is to transform the function into a smoother function with fewer local minimizers, apply an optimization algorithm to the transformed function, and trace the minimizers back to the original function. A transformed function is a coarse approximation to the original function, with small and narrow minimizers being removed, while the overall structure of the function is maintained. This property allows the optimization algorithm to skip less interesting local minimizers

and to concentrate on regions with average low-function values where a global minimizer is most likely to be located.

The smoothing transform, called the Gaussian transform, depends on a parameter  $\lambda$  that controls the degree of smoothing. The original function is obtained if  $\lambda = 0$ , while smoother functions are obtained as  $\lambda$  increases. The Gaussian transform  $\langle f \rangle_\lambda$  of a function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is

$$\langle f \rangle_\lambda(x) = \frac{1}{\pi^{n/2} \lambda^n} \int_{\mathbb{R}^n} f(y) \exp\left(-\frac{\|y - x\|^2}{\lambda^2}\right) dy. \quad (2.1)$$

The value  $\langle f \rangle_\lambda(x)$  is an average of  $f$  in a neighborhood of  $x$ , with the relative size of this neighborhood controlled by the parameter  $\lambda$ . The size of the neighborhood decreases as  $\lambda$  decreases, so that when  $\lambda = 0$ , the neighborhood is the center  $x$ . The Gaussian transform  $\langle f \rangle_\lambda$  can also be viewed as the convolution of  $f$  with the Gaussian density function.

The Gaussian transform is a linear, isotone (order-preserving) operator that reduces the high-frequency components of  $f$ . Moreover, the Gaussian transform commutes with differentiation so that the Gaussian transform of the gradient (Hessian) is the gradient (Hessian) of the Gaussian transform. These properties of the Gaussian transform are not usually shared by other approaches to smoothing. For additional discussion of these properties, see Wu [31] and Moré and Wu [23, 24].

We illustrate the transformation process in Figure 2.1 with a function that is the sum of four Gaussians. The original function ( $\lambda = 0$ ) is on the left while  $\lambda = 0.3$  is on the right. Note that the original function has four maximizers but that two of these maximizers have disappeared at  $\lambda = 0.3$ , and another minimizer is likely to disappear if  $\lambda$  is increased further. Figure 2.1 shows that the original function is gradually transformed into a smoother function with fewer local maximizers and that the smoothing increases as  $\lambda$  increases.

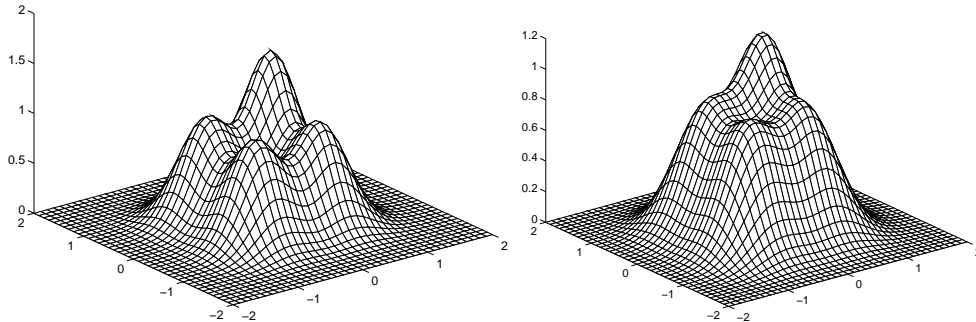


Figure 2.1: The Gaussian transform of a function. The original function ( $\lambda = 0$ ) is on the left, while  $\lambda = 0.3$  is on the right.

### 3 Computing the Gaussian Transform of Distance Geometry Functions

Computing the Gaussian transform requires the evaluation of  $n$ -dimensional integrals, but for many functions that arise in practice, it is possible to compute the Gaussian transform explicitly in terms of one-dimensional transforms. In particular, we now show that we can compute the Gaussian transform for distance geometry functions of the form (1.2) if the potential  $p_{i,j}$  is a radial function, that is, a function of the form  $p_{i,j}(x) = h_{i,j}(\|x\|)$ .

Moré and Wu [22] showed that the Gaussian transform for the distance geometry function (1.2) can be expressed in the form

$$\langle f \rangle_\lambda(x) = \sum_{i,j \in \mathcal{S}} \frac{1}{\sqrt{2\pi} r_{i,j}} \int_{-\infty}^{+\infty} (r_{i,j} + \lambda s) h_{i,j}(r_{i,j} + \lambda s) \exp\left(-\frac{1}{2}s^2\right) ds \quad (3.1)$$

where  $r_{i,j} = \|x_i - x_j\|$ . This expression is valid for all pairwise potentials of the form  $p_{i,j}(x) = h_{i,j}(\|x\|)$ . We are interested in the case where the pairwise potential  $p_{i,j}$  is given by (1.2), so that the function  $r \mapsto h_{i,j}(r)$  is defined by

$$h_{i,j}(r) = \min^2 \left\{ \frac{r^2 - l_{i,j}^2}{l_{i,j}^2}, 0 \right\} + \max^2 \left\{ \frac{r^2 - u_{i,j}^2}{u_{i,j}^2}, 0 \right\}. \quad (3.2)$$

The one-dimensional integrals that appear in (3.1) can be evaluated explicitly in special cases. In particular, when  $l_{i,j} = u_{i,j}$  for all  $(i,j) \in \mathcal{S}$ , the Gaussian transform can be expressed [23] in the form

$$\langle f \rangle_\lambda(x) = \sum_{i,j \in \mathcal{S}} \left[ (\|x_i - x_j\|^2 - \delta_{i,j}^2)^2 + 10\lambda^2 \|x_i - x_j\|^2 \right] + \gamma,$$

where  $\gamma$  is a constant that depends on  $\lambda$ .

An interesting property of the Gaussian transform is that the Gaussian transform (3.1) is infinitely differentiable whenever  $\lambda > 0$ . On the other hand, the original potential (1.2), with the pairwise potential  $p_{i,j}$  defined by (1.3), is only piecewise twice differentiable. Moré and Wu [23, 24] provide additional information on the properties of the Gaussian transform.

The one-dimensional integrals that appears in the Gaussian transform (3.1) can be approximated by Gaussian quadratures. If we use a Gaussian quadrature with  $q$  nodes, we obtain the Gauss-Hermite approximation

$$\langle f \rangle_{\lambda,q}(x) = \sum_{i,j \in \mathcal{S}} \frac{1}{r_{i,j}} \sum_{k=1}^q w_k (r_{i,j} + \lambda s_k) h_{i,j}(r_{i,j} + \lambda s_k), \quad (3.3)$$

where  $w_k$  and  $s_k$  are standard weights and nodes for the Gaussian quadrature for integrals of the form

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} g(s) \exp\left(-\frac{1}{2}s^2\right) ds.$$

The weights and nodes can be found in the tables of Stroud and Secrest [29] or computed with the `gauss` subroutine in ORTHOPOL [6].

All of our numerical results are based on the Gauss-Hermite transform (3.3). We can gain insight into this transformation by noting that

$$\langle f \rangle_{\lambda,q}(x) = \sum_{i,j \in \mathcal{S}} [h_{i,j}]_{\lambda,q}(r_{i,j}),$$

where  $[h]_{\lambda,q}$  is a function of the distance  $r$  defined by

$$[h]_{\lambda,q}(r) = \frac{1}{r} \sum_{k=1}^q w_k(r + \lambda s_k) h(r + \lambda s_k).$$

The function  $[h]_{\lambda,q}$  agrees with the piecewise twice-differentiable function

$$h(r) = \min^2 \left\{ \frac{r^2 - l^2}{l^2}, 0 \right\} + \max^2 \left\{ \frac{r^2 - u^2}{u^2}, 0 \right\} \quad (3.4)$$

for  $\lambda = 0$ , but as  $\lambda$  increases, we obtain a smoother version of the function. This can be seen clearly in Figure 3.1, where we have plotted  $[h]_{\lambda,q}$  for  $\lambda = k/2$  with  $0 \leq k \leq 3$ ,  $q = 10$ .

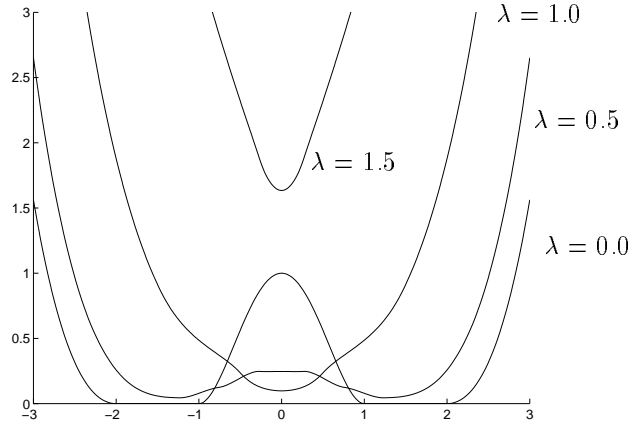


Figure 3.1: The function  $[h]_{\lambda,q}$  for  $\lambda = 0, 0.5, 1.0, 1.5$  and  $q = 10$

Figure 3.1 suggests that  $[h]_{\lambda,q}$  is convex for  $\lambda \geq \lambda_c$  for some  $\lambda_c > 0$ . This is not true in general, but holds when  $h$  is defined by (3.4). The value of  $\lambda_c$  depends on the bounds  $l$  and  $u$ , but we do not fully understand this relationship. For  $l = u$ , it is easy to show that  $\lambda_c = u/\sqrt{5}$ . Plots of  $[h]_{\lambda,q}$  for  $l < u$  suggest that  $\lambda_c \leq \sqrt{2}u$ , and therefore

$$\lambda_c \in \left[ \frac{u}{\sqrt{5}}, \sqrt{2}u \right].$$

This is consistent with Figure 3.1, where  $l = 1$ ,  $u = 2$ , and  $[h]_{\lambda,q}$  is convex for  $\lambda = 1.5$ .



The value of  $\lambda_c$  is important because if  $[h_{i,j}]_{\lambda,q}$  is monotone and convex for  $r > 0$  and  $(i,j) \in \mathcal{S}$ , then the Gauss-Hermite transform (3.3) is convex. This is usually undesirable; a preferable strategy is to choose  $\lambda$  so that only some of the functions  $[h_{i,j}]_{\lambda,q}$  are convex. We will return to this point when we discuss numerical results.

## 4 Optimization Algorithms

The algorithm that we use to solve the distance geometry problem (1.1) searches for a global minimizer of the function defined by (1.2) and (1.3) with a continuation algorithm based on the Gauss-Hermite transform  $\langle f \rangle_{\lambda,q}$ . Given a sequence of smoothing parameters

$$\lambda_0 > \lambda_1 > \cdots > \lambda_p = 0,$$

the continuation algorithm uses a local minimization algorithm to determine a minimizer  $x_{k+1}$  of  $\langle f \rangle_{\lambda_k,q}$ . The local minimization algorithm uses the previous minimizer  $x_k$  as the starting point for the search. In this manner a sequence of minimizers  $x_1, \dots, x_{p+1}$  is generated, with  $x_{p+1}$  a minimizer of  $f$  and the candidate for the global minimizer. Algorithm **dgsol** specifies the continuation algorithm.

### Algorithm **dgsol**

Choose a random vector  $x_0 \in \mathbb{R}^{m \times 3}$ .  
**for**  $k = 0, 1, \dots, p$   
    Determine  $x_{k+1} = \text{locmin}(\langle f \rangle_{\lambda_k,q}, x_k)$ .  
**end do**

In our notation,  $\text{locmin}(\langle f \rangle_{\lambda_k,q}, x_k)$  is the minimizer generated by a local minimization algorithm with the starting point  $x_k$ . The local minimization algorithm has to be chosen with some care because  $\langle f \rangle_{\lambda,q}$  is not twice continuously differentiable. The Hessian matrix is discontinuous at points where the argument of  $h_{i,j}$  coincides with either  $l_{i,j}$  or  $u_{i,j}$ . We cannot expect to avoid these discontinuities, in particular, if  $l_{i,j}$  and  $u_{i,j}$  are close.

For **locmin** we used a limited-memory variable metric algorithm of the form

$$x_{k+1} = x_k - \alpha_k H_k \nabla f(x_k),$$

where  $\alpha_k > 0$  is the search parameter, and the approximation  $H_k$  to the inverse Hessian matrix is stored in a compact representation that requires the storage of only  $2n_v$  vectors, where  $n_v$  is chosen by the user. The compact representation of  $H_k$  permits the efficient computation of  $H_k \nabla f(x_k)$  in  $(8n_v + 1)n$  flops, where  $n = 3m$  is the number of variables; all other operations in an iteration of the algorithm require  $11n$  flops.

We used the variable-metric limited-memory code **vmlm** in **MINPACK-2**. For additional information on this code, see

<http://www.mcs.anl.gov/home/more/minpack2>

The performance of the **vlm** code depends on the amount of memory specified by  $n_v$  and on the tolerances  $\tau_r$  and  $\tau_a$ . We used  $n_v = 10$ . The tolerances  $\tau_r$  and  $\tau_a$  specify the accuracy of the minimizer; **vlm** terminates with an iterate  $x$  if the code decides that either the relative convergence test

$$|f(x) - f(x^*)| \leq \tau_r |f(x^*)|$$

or the absolute convergence test

$$\max\{|f(x)|, |f(x^*)|\} \leq \tau_a$$

is satisfied for some minimizer  $x^*$  of  $f$ . In our numerical results we used  $\tau_r = \tau_a = 10^{-8}$ , which are not considered stringent values.

The random vector  $x_0 \in \mathbb{R}^{m \times 3}$  used for algorithm **dgsol** depends on the distance data. In particular, we chose the coordinates  $x_i \in \mathbb{R}^3$  of the starting point so that  $\|x_i - x_j\| = \delta_{i,j}$  for some  $(i, j)$  in  $\mathcal{S}$ . Algorithm **struct** specifies the starting point.

**Algorithm struct**

Set  $\mathcal{L} = \{1, \dots, m\}$ .

**do until**  $\mathcal{L}$  is empty

    Choose  $i \in \mathcal{L}$ .

    Set  $\mathcal{M}_i = \{j : (i, j) \in \mathcal{S}, j \in \mathcal{L}\}$ .

    For each  $j \in \mathcal{M}_i$ , generate  $x_j \in \mathbb{R}^3$  such that  $\|x_i - x_j\| = \delta_{i,j}$ .

    Remove  $i$  from  $\mathcal{L}$ .

**end do**

The starting point generated by this algorithm satisfies at least  $m - 1$  distance constraints, where  $m$  is the number of atoms. Thus, the starting point is a solution to the distance geometry problem if  $\mathcal{S}$  contains less than  $m$  constraints.

We also experimented with starting points that were chosen randomly, but since our results were not strongly dependent on the method used to generate the starting points, we present results for only the method specified above.

## 5 Computational Experiments

In our computational experiments we studied the distance geometry problem (1.1) with the pairwise potential  $p_{i,j}$  defined by (1.3). We used the **dgsol** algorithm as outlined in Section 4 and the Gauss-Hermite transform (3.3) with  $q = 10$  nodes in the Gaussian quadrature.

We tested **dgsol** on data derived from protein fragments of a DNA-binding protein [10, 27] available (ID code 1GPV) in the PDB data bank. We considered protein fragments with 100 and 200 atoms. For each fragment, we generated a set of distances  $\{\delta_{i,j}\}$  by using all distances between the atoms in the same residue as well as those in the neighboring residues.

Formally, if  $\mathcal{R}_k$  is the  $k$ th residue, then

$$\mathcal{S} = \{(i, j) : x_i \in \mathcal{R}_k, x_j \in (\mathcal{R}_k \cup \mathcal{R}_{k+1})\} \quad (5.1)$$

specifies the set of distances. This is not the only way to generate the sparse set  $\mathcal{S}$ . For example, Le Grand, Elofsson, and Eisenberg [20] generate  $\mathcal{S}$  by setting

$$\mathcal{S} = \{(i, j) : \|x_i - x_j\| \leq c\} \quad (5.2)$$

for some cutoff  $c > 0$ .

The main aim of the computational experiments is to show that the **dgsol** code, which is based on Gaussian smoothing, provides a reliable and efficient approach to the solution of the distance geometry problem (1.1). In our computational results, a set of coordinates  $x \in \mathbb{R}^{m \times 3}$  solves the distance geometry problem (1.1) if

$$(1 - \tau_d)l_{i,j} \leq \|x_i - x_j\| \leq u_{i,j}(1 + \tau_d), \quad (i, j) \in \mathcal{S}, \quad (5.3)$$

for some tolerance  $\tau_d$ . We used  $\tau_d = 10^{-2}$  since this tolerance reflects the accuracy available for bonds lengths [5].

A secondary aim of the computational experiments is to study the dependence of the solution structures on variations on the bounds  $l_{i,j}$  and  $u_{i,j}$ . For this reason we formulated distance geometry problem with lower and upper bounds  $l_{i,j}$  and  $u_{i,j}$  by setting

$$l_{i,j} = (1 - \varepsilon)\delta_{i,j}, \quad u_{i,j} = (1 + \varepsilon)\delta_{i,j}, \quad (5.4)$$

for some  $\varepsilon \in (0, 1)$ . With this formulation, we are able to study the behavior of the structures as  $\varepsilon$  varies over  $(0, 1)$ . We varied  $\varepsilon$  over  $[0.04, 0.16]$  since this translates into a 4 – 16% deviation from the expected value for the bond length. These variations seem to be typical [5].

In many of our numerical results we examine the performance of **dgsol** as  $\varepsilon$  and  $\lambda$  vary. In **dgsol** we use uniformly spaced smoothing parameters

$$\lambda_k = \lambda_0 \left(1 - \frac{k}{p}\right), \quad 0 \leq k \leq p.$$

The number  $p$  of continuation steps was set to

$$p = \lceil 20\lambda_0 \rceil.$$

This choice implies that the separation  $\lambda_{k+1} - \lambda_k$  between consecutive smoothing parameters is about 0.05.

The choice of  $\lambda_0$  is important. If we start with  $\lambda_0$  large, then all the information in the function is destroyed, and it is difficult to trace multiple paths. If we choose  $\lambda_0$  small

then  $\langle f \rangle_{\lambda_0, q}$  will have many minimizers. Choosing  $\lambda_0$  so that  $\langle f \rangle_{\lambda_0, q}$  has a few minimizers allows us to trace multiple paths, and thus increases the chances of determining a global minimizer.

A reasonable  $\lambda_0$  is obtained if half of the  $[h_{i,j}]_{\lambda_0, q}$  are not convex. This provides an automatic choice for  $\lambda_0$  that is not large and that works well. We can determine  $\lambda_0$  by recalling that (see Section 3) for each function  $h_{i,j}$  there is a  $\lambda_{i,j}$  such that  $[h_{i,j}]_{\lambda, q}$  is convex for  $\lambda \geq \lambda_{i,j}$ . We use

$$\lambda_{i,j} = \left( \frac{1}{\sqrt{5}} \rho_{i,j} + \sqrt{2}(1 - \rho_{i,j}) \right) u_{i,j}, \quad \rho_{i,j} = \frac{l_{i,j}}{u_{i,j}},$$

which specifies that  $\lambda_{i,j}$  is a convex combination of  $1/\sqrt{5}$  and  $\sqrt{2}$ . If  $l_{i,j} = u_{i,j}$  then  $\lambda_{i,j} = 1/\sqrt{5}$ , which we know guarantees convexity of  $[h_{i,j}]_{\lambda, q}$ . This observation is important because in our data  $l_{i,j} \approx u_{i,j}$ . We have verified, by plots of  $[h_{i,j}]_{\lambda, q}$  similar to those in Figure 3.1, that for this choice of  $\lambda_{i,j}$ , the function  $[h_{i,j}]_{\lambda, q}$  is convex for  $\lambda \geq \lambda_{i,j}$ . It would be interesting to obtain a formal proof of this result.

Given  $\lambda_{i,j}$  as defined above, we now choose  $\lambda_0$  as the median of all the  $\lambda_{i,j}$ . With this choice, half of the pairwise functions  $[h_{i,j}]_{\lambda_0, q}$  should not be convex. Hence, the initial function  $\langle f \rangle_{\lambda_0, q}$  is smooth but not necessarily convex.

## Experiment 1

In our first computational experiment we compare **dgsol** with **vmlm** from a set of 100 random starting points generated by algorithm **struct** of Section 4. We did this comparison because multi-starts with a local optimization code is a standard approach to solving global optimization problems. Comparisons with simulated annealing and genetic algorithms would also be of interest but are unlikely to perform better than multi-starts unless they also rely on optimization software to produce accurate structures.

We conducted two tests with  $\varepsilon = 0.04$ , one with **vmlm** and the other with **dgsol**. We compare the quality of the solutions obtained by **vmlm** and **dgsol** by computing the potential function (1.2) at the final iterate of the algorithm. These function values are then sorted and plotted in Figure 5.1.

An immediate observation that can be made from Figure 5.1 is that the potential function (1.2) has at least 100 distinct minimizers. We justify this observation by noting that all the minimizers obtained by the **vmlm** algorithm have distinct function value. This observation is of interest because it is usually difficult to find the global minimizer when the optimization problem has many minimizers.

The results in Figure 5.1 show that the **vmlm** algorithm fails to find the global minimizer in all cases. This is perhaps not surprising because the **vmlm** code is a local minimization algorithm. Nevertheless, we expected to find the global solution in at least a few cases.

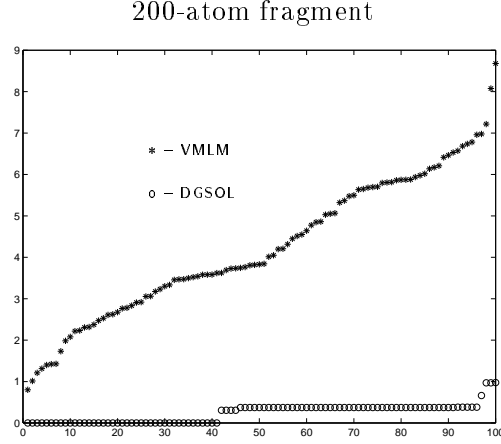


Figure 5.1: Potential function values for multi-start **vmlm** and **dgsol** for  $\varepsilon = 0.04$

However, the results in Figure 5.1 show that **vmlm** is able to find only local minimizers with relatively high function values; in all cases the potential function value is at least 0.5.

The results in Figure 5.1 also show that the smoothing approach of **dgsol** works quite well for this problem and is able to find the global solution in 41 cases. Also note that in all cases **dgsol** finds a global minimizer or a local minimizer with low function value.

## Experiment 2

In our second experiment we compare the performance of the multi-start **vmlm** with **dgsol** for problems with  $\varepsilon > 0$  and for both the 100-atom and 200-atom fragments. In each case we used the 100 random starting points generated by algorithm **struct** of Section 4 and counted the number of (global) solutions found by each algorithm. Recall that for these results we count a set of coordinates  $x \in \mathbb{R}^{m \times 3}$  as a solution to the distance geometry problem (1.1) if (5.3) is satisfied with  $\tau_d = 10^{-2}$ . Results for this experiment appear in Table 5.1.

Table 5.1: Distance geometry solutions obtained by **vmlm** and **dgsol**

100-atom fragment		
$\varepsilon$	<b>vmlm</b>	<b>dgsol</b>
0.04	1	80
0.08	1	74
0.12	8	100
0.16	47	100

200-atom fragment		
$\varepsilon$	<b>vmlm</b>	<b>dgsol</b>
0.04	0	41
0.08	0	66
0.12	2	97
0.16	9	100

The results in Table 5.1 show that **dgsol** is significantly more reliable than the multi-started **vmlm** for both the 100-atom fragment and the 200-atom fragment. For both algorithms the reliability increases with  $\varepsilon$ . This result is to be expected because as  $\varepsilon$  increases, the measure of the solution set also increases. In other words, if  $x \in \mathbb{R}^n$  satisfies (1.1) for  $l_{i,j}$  and  $u_{i,j}$  specified by (5.4), then  $x$  also satisfies (1.1) for all larger  $\varepsilon$ .

Note that the reliability of both algorithms decreases as we go from the 100-atom fragment to the 200-atom fragment. This result is to be expected because the number of minimizers of the distance geometry problem also increases as the number of atoms increases.

We emphasize that we have been using **dgsol** with 100 starting points to test the reliability of **dgsol**. In practice we can expect to find a global minimizer after at most *six* starting points. This rule of thumb is justified by the results in Table 5.1, which show that in all cases we have 40% reliability, and thus a standard calculation shows that after six trials we have a 95% chance of finding a global minimum.

### Experiment 3

In general, the distance geometry problem (1.1) can have many solutions, so there is no reason to expect that the structures generated by **dgsol** will agree with the structure that was used to generate the data. In this experiment we study the relationship between the structures obtained for various  $\varepsilon$  and the original data.

We compare structures by measuring the deviation between the coordinates and the distances for the generated structure and the original structure. A standard measure for comparing structures is the coordinate RMSD (root-mean-square-deviation)

$$E_C = \min \left\{ \left( \frac{1}{m} \sum_{i=1}^m \|y_i - Qx_i\|^2 \right)^{1/2} : Q \in \mathbb{R}^{3 \times 3}, \text{ orthogonal} \right\}, \quad (5.5)$$

where  $m$  is the number of atoms in the structure. Optimal superposition by translation is assured if the structures  $\{x_i\}$  and  $\{y_i\}$  are translated so their center of gravity is at the origin. In Table 5.2 we present the results of computing  $E_C$  for the global solutions found by **dgsol**.

The computation of the coordinate error  $E_C$  is known as the orthogonal Procrustes problem in the numerical analysis literature;  $E_C$  can be computed accurately and efficiently from the singular value decomposition of the  $3 \times 3$  matrix  $X^T Y$ , where  $X = [x_1, \dots, x_m]$  and  $Y = [y_1, \dots, y_m]$ . For details see, for example, Golub and VanLoan [9, page 582].

The coordinate error  $E_C$  is commonly used to measure the deviation between structures. In particular, many researchers require that structures have an  $E_C$  of 1-2 Å to be considered similar, while others only require an  $E_C$  of 2-3 Å. These criteria are not universally accepted

Table 5.2: Coordinate error  $E_C$  for 100-atom (left) and 200-atom (right) fragments

	$E_C$ (RMSD)			$E_C$ (RMSD)	
$\varepsilon$	Min	Ave	$\varepsilon$	Min	Ave
0.04	0.063	0.067	0.04	1.5	1.7
0.08	0.11	0.12	0.08	1.5	1.9
0.12	0.27	0.60	0.12	1.4	2.2
0.16	0.37	1.0	0.16	0.7	2.9

since  $E_C$  has a number of deficiencies. In particular,  $E_C$  is dependent on the scaling of the coordinates. For a discussion of these deficiencies, see Mairov and Crippen [21].

If we accept the view that proteins with  $E_C$  of 2-3 Å are similar, then the results in Table 5.2 show that, on the average, **dgsol** is able to find structures that are similar to the original structure. If we adopt the more stringent criterion that structures with  $E_C$  of 1-2 Å are similar, then our results show that **dgsol** finds structures that are similar if  $\varepsilon \leq 0.08$ , that is, if the lower and upper bounds differ by about 16%. If we increase  $\varepsilon$  past 0.08 then the average of  $E_C$  becomes larger than 2 Å, but, as shown by the smallest  $E_C$ , we are still able to find similar structures.

We did not expect to find small values for  $E_C$  since our data does not include all the distances, but only the distances between successive residues in the sequence. Moreover, note that we are not including all the distances within a given cutoff, as when the sparsity set  $\mathcal{S}$  is specified by (5.2).

#### Experiment 4

In the last experiment we did not consider the performance of **dgsol**. Instead, we wanted to verify, computationally, that the number of minimizers of the Gauss-Hermite transform  $\langle f \rangle_{\lambda,q}$  decreases as  $\lambda$  increases. This experiment is interesting from a theoretical viewpoint because it provides insight into the smoothing approach. We used **vm1m** with the 100 random starting points generated by algorithm **struct** on the 200-atom fragment.

The number of distinct minimizers found by **vm1m** is plotted in Figure 5.2. For these results, minimizers  $x_1$  and  $x_2$  of  $\langle f \rangle_{\lambda,q}$  are declared to be the same if

$$|\langle f \rangle_{\lambda,q}(x_1) - \langle f \rangle_{\lambda,q}(x_2)| \leq \tau_r \max \{ |\langle f \rangle_{\lambda,q}(x_1)|, |\langle f \rangle_{\lambda,q}(x_2)| \},$$

where  $\tau_r = 10^{-6}$ , or if

$$\max \{ |\langle f \rangle_{\lambda,q}(x_1)|, |\langle f \rangle_{\lambda,q}(x_2)| \} \leq \tau_a,$$

where  $\tau_a = 10^{-2}$ . In other words, the minimizers are declared to be equal if they are smaller than  $\tau_a$ , or if they are larger than  $\tau_a$  and their relative error is at most  $\tau_r$ .

The number of minimizers is sensitive to the choice of  $\tau_r$  and  $\tau_a$ , but the general trend is clear. The results in Figure 5.2 show that, as predicted by the theory, the number of minimizers of  $\langle f \rangle_{\lambda,q}$  decreases as  $\lambda$  increases. Also note that the initial drop in the number of minima is dramatic as  $\lambda$  varies in  $(0, 1)$ .

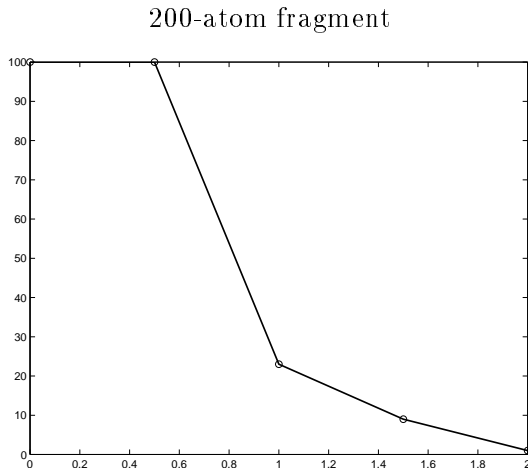


Figure 5.2: Number of minimizers of  $\langle f \rangle_{\lambda,q}$  as a function of  $\lambda$  for  $\varepsilon = 0.04$

## 6 Concluding Remarks

Our computational results suggest that protein structures can be determined by solving a distance geometry problem with **dgsol** and that the approach based on **dgsol** is significantly more reliable and efficient than multi-starts with an optimization code.

Our results also raise a number of interesting issues that we plan to address in future work. In particular, we wish to expand our testing to larger protein fragments (possibly a complete protein) and to distance data generated from NMR experiments.

Another interesting issue is the dependence of the structures on the distance data. From a mathematical viewpoint, we do not know when structures can be determined uniquely with exact, but incomplete distance data. For some results in this direction, see Hendrickson [13, 14].

## Acknowledgments

Our work has been influenced, in particular, by conversations with Paul Bash, Gordon Crippen, and Teresa Head-Gordon. Gail Pieper, as usual, deserves special thanks for her comments on the manuscript.



## References

- [1] A. ASZÓDI, M. J. GRADWELL, AND W. R. TAYLOR, *Protein fold determination using a small number of distance restraints*, in Protein Folds, A Distance Based Approach, H. Bohr and S. Brunak, eds., CRC Press, Inc., 1996, pp. 85–97.
- [2] J. M. BLANEY AND J. S. DIXON, *Distance geometry in molecular modeling*, in Reviews in Computational Chemistry, K. B. Lipkowitz and D. B. Boyd, eds., vol. 5, VCH Publishers, 1994, pp. 299–335.
- [3] A. T. BRÜNGER AND M. NILGES, *Computational challenges for macromolecular structure determination by X-ray crystallography and solution NMR-spectroscopy*, Q. Rev. Biophys., 26 (1993), pp. 49–125.
- [4] G. M. CRIPPEN AND T. F. HAVEL, *Distance Geometry and Molecular Conformation*, John Wiley & Sons, 1988.
- [5] R. A. ENGH AND R. HUBER, *Accurate bond and angle parameters for X-ray protein structure refinement*, Acta Cryst., 47 (1991), pp. 392–400.
- [6] W. GAUTSCHI, *Algorithm 726: ORTHOPOL – A package of routines for generating orthogonal polynomials and Gauss-type quadrature rules*, ACM Trans. Math. Software, 20 (1994), pp. 21–62.
- [7] W. GLUNT, T. L. HAYDEN, AND M. RAYDAN, *Molecular conformation from distance matrices*, J. Comp. Chem., 14 (1993), pp. 114–120.
- [8] ———, *Preconditioners for distance matrix algorithms*, J. Comp. Chem., 15 (1994), pp. 227–232.
- [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press, 1989.
- [10] Y. GUAN, H. ZHANG, R. N. H. KONINGS, C. W. HILBERS, T. C. TERWILLIGER, AND A. H.-J. WANG, *Crystal structure of Y41H and Y41F mutants of gene V suggest possible protein-protein interactions in the GVP-SSDNA complex*, Biochemistry, 33 (1994), p. 7768.
- [11] T. F. HAVEL, *An evaluation of computational strategies for use in the determination of protein structure from distance geometry constraints obtained by nuclear magnetic resonance*, Prog. Biophys. Mol. Biol., 56 (1991), pp. 43–78.
- [12] ———, *Distance geometry*, in Encyclopedia of Nuclear Magnetic Resonance, D. M. Grant and R. K. Harris, eds., John Wiley & Sons, 1995, pp. 1701–1710.

- [13] B. A. HENDRICKSON, *The molecule problem: Determining conformation from pairwise distances*, PhD thesis, Cornell University, Ithaca, New York, 1991.
- [14] ———, *The molecule problem: Exploiting structure in global optimization*, SIAM J. Optimization, 5 (1995), pp. 835–857.
- [15] J. C. HOCH AND A. S. STERN, *A method for determining overall protein fold from nmr distance restraints*, J. Biomolecular NMR, 2 (1992), pp. 535–543.
- [16] J. KOSTROWICKI AND L. PIELA, *Diffusion equation method of global minimization: Performance for standard functions*, J. Optim. Theory Appl., 69 (1991), pp. 269–284.
- [17] J. KOSTROWICKI, L. PIELA, B. J. CHERAYIL, AND H. A. SCHERAGA, *Performance of the diffusion equation method in searches for optimum structures of clusters of Lennard-Jones atoms*, J. Phys. Chem., 95 (1991), pp. 4113–4119.
- [18] J. KOSTROWICKI AND H. A. SCHERAGA, *Application of the diffusion equation method for global optimization to oligopeptides*, J. Phys. Chem., 96 (1992), pp. 7442–7449.
- [19] I. D. KUNTZ, J. F. THOMASON, AND C. M. OSHIRO, *Distance geometry*, in Methods in Enzymology, N. J. Oppenheimer and T. L. James, eds., vol. 177, Academic Press, 1993, pp. 159–204.
- [20] S. LE GRAND, A. ELOFSSON, AND D. EISENBERG, *The effect of distance-cutoff on the performance of the distance matrix error when used as a potential function to drive conformational search*, in Protein Folds, A Distance Based Approach, H. Bohr and S. Brunak, eds., CRC Press, Inc., 1996, pp. 105–113.
- [21] V. N. MAIROV AND G. M. CRIPPEN, *Size-independent comparison of protein three-dimensional structures*, Proteins: Struct. Func. Genetics, 22 (1995), pp. 273–283.
- [22] J. J. MORÉ AND Z. WU,  *$\varepsilon$ -optimal solutions to distance geometry problems via global continuation*, in Global Minimization of Nonconvex Energy Functions: Molecular Conformation and Protein Folding, P. M. Pardalos, D. Shalloway, and G. Xue, eds., American Mathematical Society, 1995, pp. 151–168.
- [23] ———, *Global continuation for distance geometry problems*, Preprint MCS-P505-0395, Argonne National Laboratory, Argonne, Illinois, 1995. Accepted for publication in the SIAM Journal on Optimization.
- [24] ———, *Smoothing techniques for macromolecular global optimization*, in Nonlinear Optimization and Applications, G. D. Pillo and F. Giannessi, eds., Plenum Press, 1996, pp. 297–312.

- [25] L. PIELA, J. KOSTROWICKI, AND H. A. SCHERAGA, *The multiple-minima problem in the conformational analysis of molecules: Deformation of the protein energy hypersurface by the diffusion equation method*, J. Phys. Chem., 93 (1989), pp. 3339–3346.
- [26] H. A. SCHERAGA, *Predicting three-dimensional structures of oligopeptides*, in Reviews in Computational Chemistry, K. B. Lipkowitz and D. B. Boyd, eds., vol. 3, VCH Publishers, 1992, pp. 73–142.
- [27] M. M. SKINNER, H. ZHANG, D. H. LESCHNITZER, Y. GUAN, H. BELLAMY, R. M. SWEET, C. W. GRAY, R. N. H. KONINGS, A. H.-J. WANG, AND T. C. TERWILLIGER, *Structure of the gene V protein of bacteriophage F1 determined by multi-wavelength X-ray diffraction on the selenomethionyl protein*, Proc. Nat. Acad. Sci. USA, 91 (1994), p. 2071.
- [28] M. J. SMITH-BROWN, D. KOMINOS, AND R. M. LEVY, *Global folding of proteins using a limited number of distance constraints*, Protein Engng., 6 (1993), pp. 605–614.
- [29] A. H. STROUD AND D. SECREST, *Gaussian Quadrature Formulas*, Prentice-Hall, Inc., 1966.
- [30] A. E. TORDA AND W. F. VAN GUNSTEREN, *Molecular modeling using nuclear magnetic resonance data*, in Reviews in Computational Chemistry, K. B. Lipkowitz and D. B. Boyd, eds., vol. 3, VCH Publishers, 1992, pp. 143–172.
- [31] Z. WU, *The effective energy transformation scheme as a special continuation approach to global optimization with application to molecular conformation*, SIAM J. Optimization, 6 (1996), pp. 748–768.